

Enhancing Text-to-SQL Parsing through Question Rewriting and Execution-Guided Refinement

Wenxin Mao¹, Ruiqi Wang¹, Jiyu Guo¹, Jichuan Zeng³, Cuiyun Gao^{1,2}, Peiyi Han^{1,2}, Chuanyi Liu^{1,2*}

¹Harbin Institute of Technology (Shenzhen)

²Peng Cheng Laboratory

³The Chinese University of Hong Kong

Correspondence: {23s051015,200111606,220110126}@stu.hit.edu.cn, jczeng@cse.cuhk.edu.hk,
gaocuiyun@hit.edu.cn, hanpeiyi@hit.edu.cn, liuchuanyi@hit.edu.cn

Abstract

Large Language Model (LLM)-based approach has become the mainstream for Text-to-SQL task and achieves remarkable performance. In this paper, we augment the existing prompt engineering methods by exploiting the database content and execution feedback. Specifically, we introduce DART-SQL, which comprises two key components: (1) Question Rewriting: DART-SQL rewrites natural language questions by leveraging database content information to eliminate ambiguity. (2) Execution-Guided Refinement: DART-SQL incorporates database content information and utilizes the execution results of the generated SQL to iteratively refine the SQL. We apply this framework to the two LLM-based approaches (DAIL-SQL and C3) and test it on four widely used benchmarks (Spider-dev, Spider-test, Realistic and DK). Experiments show that our framework for DAIL-SQL and C3 achieves an average improvement of 12.41% and 5.38%, respectively, in terms of the execution accuracy (EX) metric.

1 Introduction

Text-to-SQL is crafted to convert natural language questions into executable SQL queries, creating a user-friendly interface for relational databases that enhances both usability and query efficiency (Wang et al., 2022). It advances multiple facets of data management, including database accessibility and web design flexibility.

Recent researches (Rajkumar et al., 2022; Liu and Tan, 2023; Nan et al., 2023; Tai et al., 2023; Sun et al., 2023) prominently focus on the utilization of emerging LLMs, employing zero-shot and few-shot learning techniques to harness the knowledge and generalization capabilities of the LLM. The DIN-SQL approach, introduced by Pourreza and Rafiei (2023), breaks down intricate Text-to-SQL task into smaller subtasks, enhancing the performance of the LLM in the reasoning process.

When combined with GPT-4, DIN-SQL achieves an impressive score of 85.3 on the Spider-test dataset. Building upon the foundations laid by DIN-SQL, Dong et al. (2023) develop C3, facilitated by Clear Prompting (CP), Calibration with Hints (CH), and Consistent Output (CO). C3 achieves equivalent performance to DIN-SQL through Zero-Shot Learning. Further advancements in this domain are exemplified by the DAIL-SQL approach proposed by Gao et al. (2023). DAIL-SQL adopts a comprehensive approach, constructing prompts from three distinct perspectives: problem description, sample selection, and sample description. Guided by clear mission interpretation, high relevance, and a strategic SC (strategy and clarification) strategy, DAIL-SQL achieves the noteworthy feat of securing the second-highest score on the Spider Leaderboard.

However, most popular benchmarks like Spider (Yu et al., 2019) and WikiSQL (Zhong et al., 2017) focus on the database schema, with only a few rows of data in the database (Li et al., 2023c). There exists a significant gap between academic research and real-world applications, resulting in existing approaches that place more emphasis on database schema over database content information. The inherent fuzziness in natural language also poses a challenge when linking the entity names referenced in the question and the corresponding data in the database. For instance, when addressing specific conditions like the representation of gender, database entries may employ diverse representation such as ‘0’, ‘f’ and others, whereas natural language questions often use the term ‘female’. This inherent fuzziness introduces uncertainty that can significantly impact the execution efficiency of the generated SQL statements. Simultaneously, it is noteworthy that existing methods often overlook the valuable feedback generated after the execution of SQL statements. There is a lot of scope to exploit the feedback to further fine-tune the generated SQL statements.

*Corresponding author

In this work, we present a unified framework, called DART-SQL¹, which is designed to mine database content information, eradicate ambiguity in natural language questions and iteratively enhance Text-to-SQL performance by incorporating execution information. In detail, DART-SQL leverages the power of database content information through a multifaceted approach. Initially, it rewrites natural language questions by seamlessly integrating database content information, ensuring the revised questions align with the database and eliminate ambiguity. Additionally, DART-SQL introduces database content data into prompts, thereby enhancing Text-to-SQL performance. Finally, DART-SQL employs an AI agent to evaluate and analyze the results of SQL execution, progressively refining the SQL statement with the feedback from the AI agent. Through this iterative correction process, we achieve the self-training effect, enhancing the overall efficacy of DART-SQL. Experiments show that our framework for DAIL-SQL and C3 achieves an average improvement of 12.41% and 5.38%, respectively, in terms of the EX metric. Besides, Our framework DART-SQL achieves an EX score of 82.5 when applied to DAIL-SQL, placing sixth on the Spider Leaderboard².

In summary, our main contributions are:

- We observe that both existing research and benchmarks have very few database content information, which differ from real-world application significantly .
- We introduce DART-SQL, a solution that not only tackles the three aforementioned shortcomings but also effectively enhances Text-to-SQL performance.
- Extensive experiments demonstrate that our method can be built on top of a range of existing strong methods and improve their performance on Text-to-SQL task

2 Related Work

2.1 Template Matching

In the early stages, Text-to-SQL is approached using rule-based approach (Mahmud et al., 2015).

¹DART-SQL: Database-Aware Refinement and Transformation for Text-to-SQL with Large Language Models

²The top five methods on Spider Leaderboard are based on GPT-4. Due to resource limitations, we apply GPT-3.5 in our study, which hinders the performance of our framework.

Methods in this category utilize the strong rule-based feature of the SQL to design different templates for different query problems, then select the template according to the problem, and then derive the corresponding template SQL statements.

2.2 Seq2Seq Models

Owing to the inherent alignment between Text-to-SQL task and the N -> M machine translation task, many researchers adopt Seq2Seq models to address Text-to-SQL task. They first build an encoder to understand the database table structure and the input user query problem, such as the LSTM-based method (Yu et al., 2018), and the Transformer-based method (Kelkar et al., 2020), and then use a decoder to predict the target SQL statement, like the sketch-based method (Yu et al., 2018), and the generation-based method (Huang et al., 2021).

2.3 Pre-trained Models

The utilization of pre-trained language models can augment the parsing capability of Text-to-SQL. As the textual data, tabular data and SQL queries are heterogeneous, it is non-trivial but important to develop a joint reasoning framework over the three types of data (Qin et al., 2022). Qi et al. (2022) incorporate various associations such as schema encoding, schema linking, and syntactic dependencies of a problem into a unified relational representation by constructing ternary groups. Li et al. (2023a) propose a ranking-enhanced encoding and skeleton-aware decoding framework to decouple the schema linking and the skeleton parsing . The MIGA framework, proposed by Fu et al. (2023), introduces a two-stage unified multi-task generation approach that promotes compatibility between dialogue-based Text-to-SQL methods and generative pre-trained language models. Li et al. (2023b) introduce GRAPHIX-T5, which excels in capturing relational structural information while preserving the robust contextual encoding abilities inherited from the pre-trained T5 model.

2.3.1 Integrity-aware Decoding

Integrity-aware decoding verifies the code integrity of the current sequence by various means during the decoding phase, typically beam search, such as executing the statement or constructing an AST (Abstract Syntax Tree). Sequences that do not pass the check are removed from the process. Wang et al. (2018) introduce execution-guided decoding, which discards the current sequence if it is

executable but produces an execution error or outputs nothing. [Lyu et al. \(2020\)](#) verifies the SELECT clauses and conditions separately on the current step of beam search. A clause or a condition fails the check if it corresponds to no rows in the database. SeaD ([Xu et al., 2023](#)) improves the execution-guided decoding by changing the beam length of most tokens to 1, and dropping the aggregate operators during the decoding process. PICARD ([Scholak et al., 2021](#)) offers 4 different levels of checking, ranging from checking the spelling of identifiers only, to partially constructing an AST based on the current SQL statement to verify its syntactic and semantic correctness, and rejects wrong predictions.

2.4 Large Language Models

LLMs have excellent performance and generalization ability, and their outstanding performance in zero-shot and few-shot learning has attracted the attention of the industry. [Liu et al. \(2023\)](#) provide an initial comprehensive analysis of ChatGPT’s Text-to-SQL capabilities, revealing its powerful capabilities in Text-to-SQL. Recent researches have focused on exploring the appropriate form of prompt organization to maximize the performance of large models in Text-to-SQL.

2.4.1 Table Representation

The table and column name associated with the entered user query question is an important part of the prompt in Text-to-SQL. The key to table representation is to select only the database tables and columns that will be used to generate the target SQL. The study by [Poureza and Rafiei \(2023\)](#) focuses on identifying and extracting pertinent columns and tables from the database schema based on a chain-of-thought template. Meanwhile, [Dong et al. \(2023\)](#) simplify the table representation by schema linking which recalls relevant tables and columns to filter out irrelevant tables and columns, thereby reducing the number of tokens required for the prompt and improving context clarity. Also, how database information is displayed in the prompt can affect the performance of the model. For example, [Nan et al. \(2023\)](#) find that describing table structure information as CREATE statements can yield improvements.

2.4.2 Demonstration Selection

The demonstration provided in the prompt can improve the performance of the LLMs in Text-to-

SQL. [Poesia et al. \(2022\)](#) point out that presentation retrieval based on similarity is effective in a cross-domain environment, while [Levy et al. \(2023\)](#) demonstrate that diverse demonstrations benefit in-context compositional generalization. In terms of the basis for selecting demonstration, it is mainly divided into those based on the similarity of the input question ([Guo et al., 2023](#)), those based on the output SQL ([Nan et al., 2023](#)), and those that consider both the input question and the output SQL ([Gao et al., 2023](#)).

3 Methodology

In this section, we present our method that leverages database contents to enhance the natural language questions and to provide the LLM with feedback using execution results. This enables the model to learn from its own mistakes and improve its SQL generation. Our framework can be applied on top of various existing prompt designs to improve the performance. These prompts usually consist of at least two components: the schema of the database and a few examples of natural language questions along with their corresponding SQL SELECT statements. We omit the details of these prompt designs as they are not the main focus of our study. Readers can refer to the related works on 2.4 for more information.

In a nutshell, our framework consists of the following steps: ①**Question Rewriting**: We choose several rows from the table in the corresponding database related to the question as relevant database content information to rewrite the natural language question. ②**Context-Aware Prompt Generation**: For each table, we select a few rows relevant to the question, and sample some rows from the remaining ones. Such content information is then combined with an initial prompt generated by an existing prompt design to form a new context-aware prompt. ③**Execution-Guided Refinement**: The LLM takes the finalized context-aware prompt as input, and generates an SQL query. Our framework then executes it locally and offers feedback to assist the LLM to refine the statement.

3.1 Question Rewriting

To eliminate ambiguity in natural language questions and gain deeper insights into the database, we introduce an additional stage “Question Rewriting” to the conventional prompting method, which optimizes the natural language questions by incor-

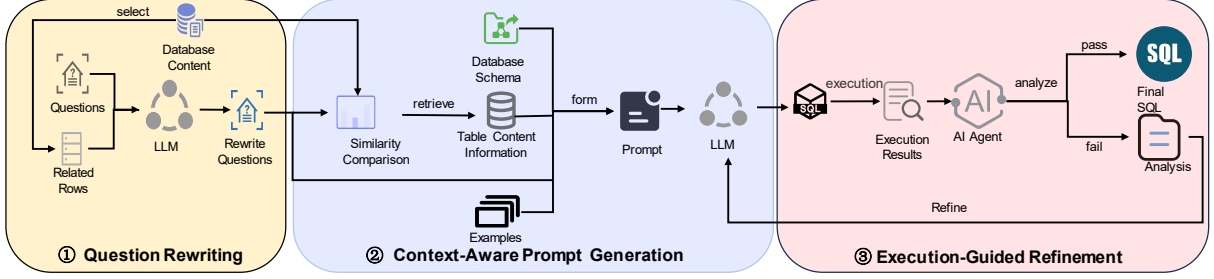


Figure 1: Overview of DART-SQL.

porating database information.

To clarify the “Question Rewriting” task for the LLM and make effective use of the database content information, we design the following prompt:

Instruction. In this section, we give a clear definition of the “Question Rewriting” task. The details are as follows: “Formulate natural language queries for database data based on the provided information. Ensure the questions you rewrite are clear, concise, and aligned with the data specifications within the database. If you think the sentence is clear enough, you can return to the original without rewriting it.”

Examples. In this section, we summarize the three specifications for the model to follow when rewriting questions and illustrate them with concrete examples. Firstly, we instruct the model on how to rewrite questions exhibiting fuzziness through the specification 1. Secondly, Specification 2 serves as a reminder for the model not to omit key information during the rewrite process. For instance, database content information like “Note that 1 stands for yes, and 0 stands for NO in the tables” in the original question might be omitted after rewriting. Finally, through Specification 3, we ensure that the model avoids making meaningless or potentially harmful modifications to questions that are already sufficiently clear.

Input. In this section, we take the initial natural language question and related database content information as input. Previous methods focused on schema and structure but neglected effective use of database content. To address this, we first identify the database corresponding to the natural language question based on annotations in the dataset. Next, we retrieve the names of all tables in the database and proceed to search for the top K records (due to prompt length limitations, we set $K = 5$) in each table. Finally, we incorporate this information into the prompt. Surprisingly, thanks to the

LLMs’ summarization and generalization capabilities, the model mostly identified the relevant table and refined the problem based on specific database statements, ensuring consistency in entity names, cases, abbreviations, and filter criteria for attribute names.

The Question Rewriting module greatly reduces problem ambiguity. Due to the limited context length of the LLM, the original Text-to-SQL task prompt includes only limited database information. By introducing the additional Question Rewriting task, we can effectively incorporate underutilized database content information into the natural language problems during rewriting. This allows for the utilization of more database information within the constraints of a limited context length.

3.2 Context-Aware Prompt Generation

After rewriting the question, the initial prompt can be generated using an existing prompt design. Our framework then refines the prompt to become the context-aware prompt by adding additional portions of the table content. Intuitively, this refinement aims to reduce value errors by providing the LLM with examples of the data format in the database. Moreover, these data samples also complement the database schema, assisting the LLM to understand the database structure.

At most N rows of each table are included in the prompt. If a table has more than N rows, half of the rows are selected based on their relevance to the question, and the other half are sampled from the remaining rows. We set $N = 15$ for our method.

To calculate the relevance of a row inside a table to a question Q_i , we first convert the row into a JSON dictionary D_i , with column names as keys and the cell contents as values. The dictionary and the question are tokenized respectively, resulting in $D_i = \{D_{i1}, \dots, D_{i|D_i|}\}$ and $Q_i = \{Q_{i1}, \dots, Q_{i|Q_i|}\}$. Next, we compute the

sentence embeddings of D_i and Q_i :

$$E(D_i) = \frac{1}{|D_i|} \sum_{k=1}^{|D_i|} \text{GloVe}(D_{ik})$$

$$E(Q_i) = \frac{1}{|Q_i|} \sum_{k=1}^{|Q_i|} \text{GloVe}(Q_{ik})$$

where $\text{GloVe}(x)$ refers to the GloVe embedding (Pennington et al., 2014)³ of token x . Finally, the relevance is calculated as the cosine similarity between $E(D_i)$ and $E(Q_i)$ before the most relevant rows as well as some random rows are integrated into the prompt.

3.3 Execution-Guided Refinement

After receiving the processed prompt, the LLM then outputs a single SELECT statement to answer the question. We then start a self-refinement procedure that executes the statement on the database and ask the LLM to analyze the result, until it thinks the results are correct or the procedure is repeated 3 times. During the refinement process, the LLM can only refer to the initial prompt, the latest generated SQL statement, and its execution results. The complete procedure is shown in Algorithm 1.

In the refinement process, the SELECT statement is executed. If the execution is successful, the results (truncated at $M = 15$ rows) are analyzed by the LLM. Otherwise, the LLM receives the exception information and outputs a new statement addressing the issue. If the exception is related to column names, including ambiguity and no-such-column errors, we additionally provide the LLM with a list of tables containing the column.

Furthermore, the LLM is forbidden from using certain keywords including LEFT JOIN and SELECT *. If the output SQL statement contains such keywords, the LLM is immediately asked to regenerate it. The rationale is that upon manually analyzing the final outputs, most statements with these keywords produce incorrect results. This suggests that the LLMs we test are not proficient at utilizing these SQL features.

³We choose GloVe for obtaining sequence embeddings since our baseline methods such as DAIL-SQL use GloVe to compute embeddings during schema linking.

Algorithm 1 Self-Refinement Algorithm

Input: initial SELECT statement SQL_0

Output: refined statement SQL

$SQL \leftarrow SQL_0$

for $i = 1$ to 3 **do**

 Reset the conversation history, retaining only the initial prompt from user and SQL from LLM

$res \leftarrow$ execution result of SQL

if res indicates a column-name error **then**

$c \leftarrow$ incorrect column name

$tables \leftarrow$ list of tables with column c

 Send $tables$ and res to LLM

else if res is other types of exception **then**

 Send res to LLM

else $\triangleright res$ stores execution results

 Send at most M rows in res to LLM

if LLM thinks res is correct **then**

return SQL

end if

end if

$SQL \leftarrow$ refined SQL by LLM

end for

return SQL

4 Experiments

4.1 Experimental Setup

Dataset. We evaluate DART-SQL on three datasets, **Spider**⁴ (Yu et al., 2019), **Spider-DK** (Gan et al., 2021), and **Spider-Realistic** (Deng et al., 2021). Spider is a large-scale cross-domain Text-to-SQL dataset. It comprises 10181 annotated questions, corresponding to 5693 distinct complex SQL queries and 200 databases with multiple tables. The training, development, and test sets of Spider include 8659, 1034, and 2147 examples and 146, 20, and 40 databases respectively. Our approach selects pertinent few-shot samples from the training set. Spider-DK and Spider-Realistic are more challenging variants of the original Spider dataset. Spider-DK features 535 instances and emphasizes the importance of domain knowledge, while Spider-Realistic consists of 508 instances from Spider with explicit mentions of column names removed.

Metric. Following previous studies, we use exact-set-match accuracy (**EM**) and execution accuracy (**EX**) to evaluate our framework. EM breaks down the predicted SQL statement and the ground-truth into two sequences of SQL clauses, and

⁴As of January 2024, the authors have publicized the test set of Spider, available at the original download link.

checks if they are identical after removing the values in the statements. EX compares the execution results of the predicted statement and the ground truth on the corresponding database. The test suite we use was introduced in (Zhong et al., 2020).

Model and hyper-parameter. We use the 0613 version of GPT-3.5-turbo-16k as our LLM via the OpenAI API for all experiments, including the baseline methods⁵. We set the number of chat completion choices n to 1 and the temperature T to 0.0. When self-consistency (Wang et al., 2023) is enabled, we set $n = 5$ and $T = 1.0$ for SQL generation instead, and leave the parameters unchanged for text generation.

Base prompt. We select DAIL-SQL and C3 to generate the initial prompt. DAIL-SQL includes the entire database schema in the prompt, and selects few-shot examples based on both the questions and the pre-generated queries. C3 uses the LLM to perform zero-shot schema linking, incorporates guidelines into the conversation history, and features execution-based self-consistency.

4.2 Overall Performance

Table 1 and 2 display the results of our framework as well as the baseline methods on the datasets. According to the results, our framework brings noticeable performance improvement to both prompt designs in terms of execution accuracy. When DAIL-SQL is used as the base prompt, DART-SQL drastically outperforms the baseline method by 15.2% in EM and 10.4% in EX on the Spider-Realistic dataset, while also providing satisfying performance boost on other datasets. When combined with C3, DART-SQL achieves 73.8% in EX, outperforming both the baseline and the separated modules in DART-SQL. This indicates that DART-SQL has superior skills in utilizing domain knowledge to solve challenging problems in the datasets. On the other hand, self-consistency only boosts the performance in EX on the test set of Spider when using DAIL-SQL, and negatively affect the EM scores on all datasets, which means enabling self-consistency is unnecessary for our framework.

4.3 Analysis of DART-SQL Components

We conduct an ablation study on the development set of Spider with DAIL-SQL as our base prompt,

⁵The baseline methods are rerun with the authors’ open-source repositories on GPT-3.5, and the results might differ from those reported in the original paper.

and record the results in Table 3. In each experiment, one or more components of DART-SQL is disabled to validate their necessity. In addition, we perform a hyper-parameter analysis on our framework in Figure 2, regarding the modules of table content selection and execution-guided refinement.

4.3.1 Example Selection

We test the performance of DART-SQL under the zero-shot scenario, i.e. all examples of question-query pairs offered by DAIL-SQL are removed. Here, DART-SQL is used as a replacement of DAIL-SQL, providing a 2.4% increase in EX, which indicates that the LLM has the ability to learn from database patterns and refine the query statement based on the execution results. However, as the training and development sets both come from the Spider dataset so they share the same design principles and distribution, the removal of few-shot samples still degrade the overall performance, not to mention the drastic decrease of 19.9% in EM scores.

As the number of examples in the prompt increases, according to Figure 2a, the performance boost depends on S , the number of samples, being insignificant when $S > 3$. This indicates that to maintain a similar level of performance, we should offer at least 3 examples, achieving 82.6% in EX.

4.3.2 Question Rewriting

In Table 3, we observe a 0.8% decline in the EM score after removing the Question Rewriting module. Simultaneously, as shown in Tables 1 and 2, the EM scores improve with the implementation of the Question Rewriting module. We attribute this to the nature of the EM metric, which measures the degree of match between the generated SQL and the Gold SQL. The Gold SQL is manually written by the dataset creators, who understand the intent of the natural language question. Therefore, the retrieval methods and filtering conditions used in the Gold SQL align with the intent of the natural language question. By clarifying the question’s intent through the Question Rewriting module, the generated SQL will more closely align with the Gold SQL. Consequently, the Question Rewriting module usually improves the EM score.

Meanwhile, in Table 3, we find that the EX score increases after removing the Question Rewriting module. Tables 1 and 2 show that Question Rewriting module results in an increase in the EX score for the majority of datasets, albeit with a decrease

Base Prompt: DAIL-SQL	Dev		Test		Realistic		DK	
	EM	EX	EM	EX	EM	EX	EM	EX
baseline	62.1	76.2	61.4	76.9	42.1	68.9	41.7	61.5
+SC	62.1	79.4	60.3	79.3	44.1	70.3	41.1	63.6
RW	<u>63.9</u>	79.9	60.3	79.3	47.8	71.5	44.7	66.5
+SC	61.7	80.9	61.2	79.6	42.3	71.7	39.4	64.1
CAP+ER	61.3	84.0	61.0	81.6	55.9	80.7	45.0	70.7
+SC	59.3	83.5	59.9	83.0	53.3	78.1	43.6	72.0
DART-SQL	62.1	83.1	62.0	81.8	57.3	79.3	44.7	71.4
+SC	59.3	83.1	59.9	82.5	53.3	79.3	41.1	70.7

Table 1: Results using DAIL-SQL as the base prompt. SC means self-consistency while RW, CAP, and ER refer to question rewriting, context-aware prompt generation, and execution-guided refinement introduced in Section 3.1, 3.2, and 3.3 respectively. For each metric on a dataset, the highest result is both bold and underlined, while the results that are at most 1 percentage lower than it are also marked bold.

Base Prompt: C3	Dev		Test		Realistic		DK	
	EM	EX	EM	EX	EM	EX	EM	EX
baseline	46.4	81.3	46.7	79.7	44.3	77.6	40.6	66.4
RW	47.3	80.9	47.2	79.5	47.4	78.5	39.6	69.0
CAP+ER	47.0	84.4	47.6	83.0	43.9	82.9	37.8	72.1
DART-SQL	46.4	83.0	46.3	82.1	45.7	81.7	38.1	73.8

Table 2: Results using C3 as the base prompt.

in a minority of datasets. Our analysis suggests that this decline in EX scores may be attributed to two issues within the rewrite module: overwriting and the omission of key information. Despite our efforts to mitigate these issues, they continue to have some negative effects.

Finally, the experimental results demonstrate that the Question Rewriting module is overall beneficial in enhancing the performance of the method. It eliminates ambiguity in questions and improves the accuracy of the generated SQL.

4.3.3 Context-Aware Prompt Generation

As can be seen in Table 3, removing the module of context-aware prompt generation decreases the EX score by 2.7% while giving a similar EM score. Without the content samples in each table, the LLM is unaware of the data formats in the database until execution time, and its understanding of the database structure is also negatively affected.

Furthermore, we design experiments to study the impact on EX scores of different N , the maximum number of rows offered to the LLM in each table in the database. According to Figure 2b, even the contents of a single row in each table makes a significant contribution of 2.0% to the EX score. As N gradually increases, the EX score also increases

slowly. Since most tables in the Spider dataset contain no more than 15 rows, we can infer that adding more rows has rather limited impact on the performance.

4.3.4 Execution-Guided Refinement

In Table 3, when the refinement module is disabled, the EX score decreases by 2.5% while the EM score maintains a similar level. If we additionally switch off the content selection module, the EX score further decreases by 0.7%. In other words, when the two modules are combined, the overall performance boost is larger than those offered by them separately. We can conclude that the LLM can better utilize the execution results to refine the query with the help of database contents.

We design two sets of experiment on M , the number of rows in the execution results and R , the maximum number of revisions by the LLM, respectively. As we can see in Figure 2c, 3 rows are barely enough for the self-refinement process, while the performance suffers greatly from diminishing returns when $M > 7$. As for the maximum number of revisions, in Figure 2d, the LLM is capable of correcting quite a few mistakes when $R = 1$, receiving another increase of 0.6% if R is increased to 3.

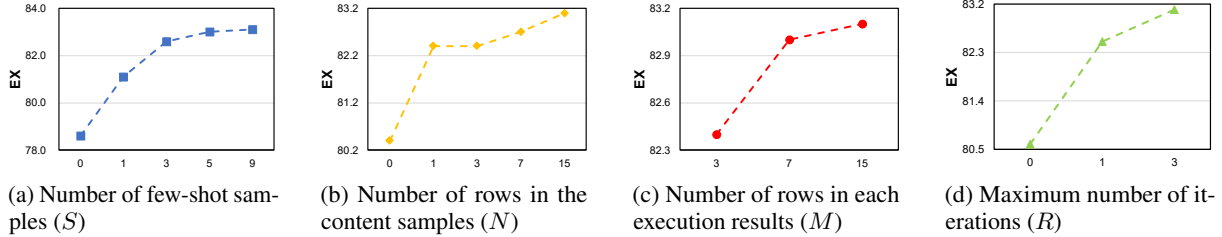


Figure 2: Hyper-parameter analysis of DART-SQL.

Question: List the name and date the battle that has lost Lettice and HMS Atalanta						
Query by DAIL: SELECT T1.name, T1.date FROM battle AS T1 JOIN ship AS T2 ON T1.id = T2.lost_in_battle WHERE T2.name = "Lettice" AND T2.ship_type = "HMS Atalanta"			✗ Execution Result: None			
Rewritten Question: List the name and date of the battle in which Lettice and HMS Atalanta were lost.						
Content Samples (Table: ship)	lost_in_battle	id	name	tonnage	ship_type	location
	8	1	<u>Lettice</u>	t	Brig	English Channel
	8	7	HMS Atalanta	225	8 gun Brig	Mid-Atlantic
Initial Query by DART: SELECT T1.name, T1.date FROM battle AS T1 JOIN ship AS T2 ON T1.id = T2.lost_in_battle WHERE T2.name = "Lettice" OR T2.name = "HMS Atalanta"			✗ Execution Result: ('Siege of Constantinople', '1235') ('Siege of Constantinople', '1235')			
Revised Query by DART: SELECT T1.name, T1.date FROM battle AS T1 JOIN ship AS T2 ON T1.id = T2.lost_in_battle WHERE T2.name = "Lettice" <u>INTERSECT SELECT T1.name, T1.date FROM battle AS T1</u> <u>JOIN ship AS T2 ON T1.id = T2.lost_in_battle</u> WHERE T2.name = "HMS Atalanta"			✓ Execution Result: ('Siege of Constantinople', '1235')			

Figure 3: Case study of DART-SQL. The modified part in the question is marked in dark red. The differences in the three SQL statements are underlined. For simplicity, we only list some relevant rows in the content samples.

Base Prompt: DAIL-SQL	Dev (EM)	Dev (EX)
DART-SQL	62.1	83.1
w/o RW	61.3	84.0
w/o CAP	62.2	80.4
w/o ER	62.4	80.6
w/o RW+CAP	61.3	80.2
w/o RW+ER	61.9	80.4
w/o CAP+ER	63.9	79.9
w/o DART-SQL (baseline)	62.1	76.2
w/o DAIL-SQL (0-shot)	42.2	78.6

Table 3: Ablation study on the dev split of the Spider dataset.

4.4 Result Analysis

Here we compare the correctness of output results of C3 as the baseline method, and DART-SQL with C3 as the base prompt design, on the Spider-DK dataset, where these methods achieve 66.4% and 73.8% execution accuracy respectively. Table 4 displays the results.

Overall, our framework addresses 66 errors but introduces 25 new mistakes compared to the baseline method. About half of the fixed errors are wrong values, since C3 provides very limited information on the database contents via schema linking,

while we directly include them in our prompt, in addition to sending feedback to the LLM with execution results.

For other errors, condition errors include bugs in the WHERE or HAVING clauses, while table/-column errors refer to mistakes in table or column names, or selecting extra columns. These bugs are much more difficult to fix than value errors, requiring deep understanding of the question and the database schema. Thanks to the question rewriting module, the clarity of the question is noticeably improved for the LLM, while the other two modules are helpful in illustrating the structure of the database. As a result, our framework also makes progress in addressing these errors.

4.5 Case Study

To explore how each module in DART-SQL encourages the LLM to generate SQL statements of higher quality, we conduct a case study on the 254th question in the Spider-Realistic dataset as shown in Figure 3.

We first generate an SQL statement with the baseline method, DAIL-SQL. The original ques-

Base Prompt: C3	Fixed	Introduced
Total	66	25
Values	34	2
Conditions	6	5
Table/column	20	16
Others	6	2

Table 4: Number of errors fixed and introduced by our framework compared to the baseline method C3, on the Spider-DK dataset.

tion is included in the prompt, and as a result of the incorrect grammar in addition to the absence of database contents, it presents difficulty for the LLM to understand the question and the database structure. With the DAIL-SQL prompt, the LLM mistakenly considers "HMS Atalanta" as a ship type, and returns a wrong query.

In the following step, we use DART-SQL to solve the problem. The question is rephrased to improve clarity, and parts of the relevant contents are integrated into the prompt. As shown in the Content Samples table in Figure 3, the two entity names are actually ship names, and the LLM generates an initial query, correctly addressing the issue in the baseline method. However, there exists one more problem that the generated statement allows duplicated rows. Once the LLM receives the results, it detects the mistake and corrects it by utilizing the INTERSECT operator, as shown in the revised query. The final execution result is correct and identical to the ground-truth output.

5 Conclusion

In this paper, we propose a novel Text-to-SQL framework, DART-SQL, which addresses the issue of insufficient utilization of database contents, ambiguity in the natural language question, and lack of self-refinement based on execution results. Our experiment on the Spider dataset and two of its variants shows significant improvement over DAIL-SQL and C3. Furthermore, we substantiated the validity and generality of our framework through a comprehensive series of experiments. Overall, our study presents a practical solution to the task of Text-to-SQL. It opens up a new research direction for those entering the field, and we anticipate more in-depth exploration in this direction in future research endeavors.

Limitations

Due to resource limitations, we are unable to conduct our research on other state-of-the-art LLMs, such as GPT-4-turbo, PaLM, and Llama, or on other base prompt designs like DIN-SQL. Our feedback to the LLM also only checks a few common issues, which can be enhanced with SQL static tools. We expect future researches to provide more helpful insights on the generated query for the LLM to further increase Text-to-SQL performance.

Ethics Statement

Our experimentation leverages the OpenAI platform to assess the efficacy of our approach, utilizing their GPT-3.5-turbo-16k model. It's worth noting that large models heavily depend on significant computing power, contributing to substantial power consumption and increased carbon dioxide emissions.

Acknowledgments

This research is supported by Natural Science Foundation of Guangdong Province (Project No. 2023A1515011959), Shenzhen-Hong Kong Jointly Funded Project (Category A, No. SGDX20230116091246007), Shenzhen Basic Research (General Project No. JCYJ20220531095214031), Shenzhen International Science and Technology Cooperation Project (No. GJHZ20220913143008015), the Major Key Project of PCL (Grant No. PCL2022A03), the National Key Research and Development Program of China under Grant 2023YFB3106504, Guangdong Provincial Key Laboratory of Novel Security Intelligence Technologies under Grant 2022B1212010005, the Major Key Project of PCL under Grant PCL2023A09, Shenzhen Science and Technology Program under Grant ZDSYS20210623091809029 and RCBS20221008093131089.

References

Xiang Deng, Ahmed Hassan Awadallah, Christopher Meek, Oleksandr Polozov, Huan Sun, and Matthew Richardson. 2021. [Structure-grounded pretraining for text-to-sql](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.

- Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, Jinshu Lin, Dongfang Lou, et al. 2023. C3: Zero-shot text-to-sql with chatgpt. *arXiv preprint arXiv:2307.07306*.
- Yingwen Fu, Wenjie Ou, Zhou Yu, and Yue Lin. 2023. Miga: a unified multi-task generation framework for conversational text-to-sql. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 12790–12798.
- Yujian Gan, Xinyun Chen, and Matthew Purver. 2021. [Exploring underexplored limitations of cross-domain text-to-sql generalization](#).
- Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2023. Text-to-sql empowered by large language models: A benchmark evaluation. *arXiv preprint arXiv:2308.15363*.
- Chunxi Guo, Zhiliang Tian, Jintao Tang, Pancheng Wang, Zhihua Wen, Kang Yang, and Ting Wang. 2023. [Prompting gpt-3.5 for text-to-sql with de-semanticization and skeleton retrieval](#).
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.
- Junyang Huang, Yongbo Wang, Yongliang Wang, Yang Dong, and Yanghua Xiao. 2021. [Relation aware semi-autoregressive semantic parsing for nl2sql](#).
- Amol Kelkar, Rohan Relan, Vaishali Bhardwaj, Saurabh Vaichal, Chandra Khatri, and Peter Relan. 2020. [Bertrand-dr: Improving text-to-sql using a discriminative re-ranker](#).
- Itay Levy, Ben Bogin, and Jonathan Berant. 2023. [Diverse demonstrations improve in-context compositional generalization](#).
- Haoyang Li, Jing Zhang, Cuiping Li, and Hong Chen. 2023a. Resdsq: Decoupling schema linking and skeleton parsing for text-to-sql. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 13067–13075.
- Jinyang Li, Binyuan Hui, Reynold Cheng, Bowen Qin, Chenhao Ma, Nan Huo, Fei Huang, Wenyu Du, Luo Si, and Yongbin Li. 2023b. Graphix-t5: Mixing pre-trained transformers with graph-aware layers for text-to-sql parsing. *arXiv preprint arXiv:2301.07507*.
- Jinyang Li, Binyuan Hui, Ge Qu, Binhua Li, Jiayi Yang, Bowen Li, Bailin Wang, Bowen Qin, Rongyu Cao, Ruiying Geng, et al. 2023c. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *arXiv preprint arXiv:2305.03111*.
- Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S. Yu. 2023. [A comprehensive evaluation of chatgpt’s zero-shot text-to-sql capability](#).
- Xiping Liu and Zhao Tan. 2023. Divide and prompt: Chain of thought prompting for text-to-sql. *arXiv preprint arXiv:2304.11556*.
- Qin Lyu, Kaushik Chakrabarti, Shobhit Hathi, Souvik Kundu, Jianwen Zhang, and Zheng Chen. 2020. [Hybrid ranking network for text-to-sql](#).
- Tanzim Mahmud, K. M. Azharul Hasan, Mahtab Ahmed, and Thwoi Hla Ching Chak. 2015. [A rule based approach for nlp based query processing](#). In *2015 2nd International Conference on Electrical Information and Communication Technologies (EICT)*, pages 78–82.
- Linyong Nan, Yilun Zhao, Weijin Zou, Narutatsu Ri, Jaesung Tae, Ellen Zhang, Arman Cohan, and Dragomir Radev. 2023. Enhancing few-shot text-to-sql capabilities of large language models: A study on prompt design strategies. *arXiv preprint arXiv:2305.12586*.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [Glove: Global vectors for word representation](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. 2022. [Synchronesh: Reliable code generation from pre-trained language models](#).
- Mohammadreza Pourreza and Davood Rafiei. 2023. Din-sql: Decomposed in-context learning of text-to-sql with self-correction. *arXiv preprint arXiv:2304.11015*.
- Jiexing Qi, Jingyao Tang, Ziwei He, Xiangpeng Wan, Yu Cheng, Chenghu Zhou, Xinbing Wang, Quanshi Zhang, and Zhouhan Lin. 2022. Rasat: Integrating relational structures into pretrained seq2seq model for text-to-sql. *arXiv preprint arXiv:2205.06983*.
- Bowen Qin, Binyuan Hui, Lihan Wang, Min Yang, Jinyang Li, Binhua Li, Ruiying Geng, Rongyu Cao, Jian Sun, Luo Si, Fei Huang, and Yongbin Li. 2022. [A survey on text-to-sql parsing: Concepts, methods, and future directions](#).
- Nitarshan Rajkumar, Raymond Li, and Dzmitry Bahdanau. 2022. Evaluating the text-to-sql capabilities of large language models. *arXiv preprint arXiv:2204.00498*.
- Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. [Picard: Parsing incrementally for constrained auto-regressive decoding from language models](#).
- Ruoxi Sun, Sercan O Arik, Hootan Nakhost, Hanjun Dai, Rajarishi Sinha, Pengcheng Yin, and Tomas Pfister. 2023. Sql-palm: Improved large language model adaptation for text-to-sql. *arXiv preprint arXiv:2306.00739*.

- Chang-You Tai, Ziru Chen, Tianshu Zhang, Xiang Deng, and Huan Sun. 2023. Exploring chain-of-thought style prompting for text-to-sql. *arXiv preprint arXiv:2305.14215*.
- Chenglong Wang, Kedar Tatwawadi, Marc Brockschmidt, Po-Sen Huang, Yi Mao, Oleksandr Polozov, and Rishabh Singh. 2018. Robust text-to-sql generation with execution-guided decoding. *arXiv preprint arXiv:1807.03100*.
- Lihan Wang, Bowen Qin, Binyuan Hui, Bowen Li, Min Yang, Bailin Wang, Binhua Li, Jian Sun, Fei Huang, Luo Si, et al. 2022. Proton: Probing schema linking information from pre-trained language models for text-to-sql parsing. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1889–1898.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models.
- Kuan Xu, Yongbo Wang, Yongliang Wang, Zujie Wen, and Yang Dong. 2023. Sead: End-to-end text-to-sql generation with schema-aware denoising.
- Tao Yu, Zifan Li, Zilin Zhang, Rui Zhang, and Dragomir Radev. 2018. Typesql: Knowledge-based type-aware neural text-to-sql generation. *arXiv preprint arXiv:1804.09769*.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2019. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task.
- Ruiqi Zhong, Tao Yu, and Dan Klein. 2020. Semantic evaluation for text-to-sql with distilled test suite. In *The 2020 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv preprint arXiv:1709.00103*.

A Prompt Examples

In the appendix, we demonstrate an example of the prompt in each step used in our framework.

A.1 Question Rewriting

/*Instructions:*/

Formulate natural language queries for database data based on the provided information. Ensure the questions you rewrite are clear, concise, and aligned with the data specifications within the database. If you think the sentence is clear enough, you can return to the original without rewriting it.

/*Here are some basic examples: */

/*Example 1: Note that rewrite question aligned with the data specifications within the database*/

/*Given the initial value from Table:*/

[Table countries: ', [(1, 'usa', 1), (2, 'germany', 2), (3, 'france', 2), (4, 'japan', 3), (5, 'italy', 2)]]

/*Rewrite the following question:What is the total number of car makers in Italy?*/

/*results:Can you provide information on the number of car manufacturers in "italy"?*/

/*Rewrite the following question:What is the count of the car models produced in the United States?*/

/*results:What is the count of the car models produced in the "usa"?*/

/*Rewrite the following question:What is the total number of car makers are there in japan?*/

/*results:How many makers who are in "japan"?*/

/*Given the initial value from Table:*/

[Table Student: ', [(1001, 'Smith', 'Linda', 18, 'F', 600, 1121, 'BAL'), (1002, 'Kim', 'Tracy', 19, 'F', 600, 7712, 'HKG'), (1003, 'Jones', 'Shiela', 21, 'F', 600, 7792, 'WAS'), (1004, 'Kumar', 'Dinesh', 20, 'M', 600, 8423, 'CHI'), (1005, 'Gompers', 'Paul', 26, 'M', 600, 1121, 'YYZ')]]

[Table Has_Pet: ', [(1001, 2001), (1002, 2002), (1002, 2003)]]

[Table Pets: ', [(2001, 'cat', 3, 12.0), (2002, 'dog', 2, 13.4), (2003, 'dog', 1, 9.3)]]

/*Rewrite the following question:Find the major and age of students who do not have a dog pet*/

/*results:Please provide the names and ages of students who do not have a pet with the pettype of "dog" in the table "Pets".*/

/*Given the initial value from Table:*/

[Table stadium: ', [(1, 'Raith Rovers', 'Štark's Park', 10104, 4812, 1294, 2106), (2, 'Ayr United', 'Somerset Park', 11998, 2363, 1057, 1477), (3, 'East Fife', 'Bayview Stadium', 2000, 1980, 533, 864)]]

[Table singer: ', [(1, 'Joe Sharp', 'Netherlands', 'You', '1992', 52, 'F'), (2, 'Timbaland', 'United States', 'Dangerous', '2008', 32, 'T'), (3, 'Justin Brown', 'France', 'Hey Oh', '2013', 29, 'T'), (4, 'Rose White', 'France', 'Sun', '2003', 41, 'F'), (5, 'John Nizininik', 'France', 'Gentleman', '2014', 43, 'T')]]

[Table concert: ', [(1, 'Auditions', 'Free choice', '1', '2014'), (2, 'Super bootcamp', 'Free choice 2', '2', '2014'), (3, 'Home Visits', 'Bleeding Love', '2', '2015')]]

[Table singer_in_concert: ', [(1, '2'), (1, '3'), (1, '5'), (2, '3'), (2, '6')]]

/*Rewrite the following question:What are the names of the singers who are not Australian citizens*/

*/*results:Which singers in the singer table have a citizenship value that is not Australia?*/*

/*Example 2: Note that you can not miss the information in original question*/

*/*Given the initial value from Table:*/*

['Table Dogs: ', [(1, 3, '1', 'ESK', 'LGE', 'Kacey', '6', '2012-01-27 05:11:53', '1', '7.57', '2017-09-08 20:10:13', '2018-03-06 16:32:11', '2018-03-25 06:58:44'), (2, 11, '0', 'BUL', 'LGE', 'Hipolito', '9', '2013-02-13 05:15:21', '0', '1.72', '2017-12-22 05:02:02', '2018-03-25 08:12:51', '2018-03-25 02:11:32'), (3, 1, '0', 'BUL', 'MED', 'Mavis', '8', '2008-05-19 15:54:49', '1', '8.04', '2017-06-25 10:14:05', '2018-03-07 21:45:43', '2018-03-25 10:25:46')]]

*/*Rewrite the following question:What are the dog name, age and weight of the dogs that were abandoned? Note that 1 stands for yes, and 0 stands for no in the tables.*/*

*/*results:What are the dog name, age and weight of the dogs in the Table Dogs that were abandoned? Note that 1 stands for yes, and 0 stands for no in the tables.*/*

/*Example 3: If you think the sentence is clear enough, you can return to the original without rewriting it*/

*/*Given the initial value from Table:*/*

['Table cars_data: ', (1, '18', 8, 307.0, '130', 3504, 12.0, 1970), (2, '15', 8, 350.0, '165', 3693, 11.5, 1970), (3, '18', 8, 318.0, '150', 3436, 11.0, 1970), (4, '16', 8, 304.0, '150', 3433, 12.0, 1970), (5, '17', 8, 302.0, '140', 3449, 10.5, 1970)]]

*/*Rewrite the following question: Among the cars that do not have the minimum horsepower , what are the make ids and names of all those with less than 4 cylinders ?*/*

*/*results:Among the cars that do not have the minimum horsepower , what are the make ids and names of all those with less than 4 cylinders ?*/*

*/*Given the initial value from Table:*/*

['Table stadium: ', [(1, 'Raith Rovers', 'Stark's Park', 10104, 4812, 1294, 2106), (2, 'Ayr United', 'Somerset Park', 11998, 2363, 1057, 1477), (3, 'East Fife', 'Bayview Stadium', 2000, 1980, 533, 864)]]
['Table singer: ', [(1, 'Joe Sharp', 'Netherlands', 'You', '1992', 52, 'F'), (2, 'Timbaland', 'United States', 'Dangerous', '2008', 32, 'T'), (3, 'Justin Brown', 'France', 'Hey Oh', '2013', 29, 'T'), (4, 'Rose White', 'France', 'Sun', '2003', 41, 'F'), (5, 'John Nizirik', 'France', 'Gentleman', '2014', 43, 'T')]]
['Table concert: ', [(1, 'Auditions', 'Free choice', '1', '2014'), (2, 'Super bootcamp', 'Free choice 2', '2', '2014'), (3, 'Home Visits', 'Bleeding Love', '2', '2015')]]
['Table singer_in_concert: ', [(1, '2'), (1, '3'), (1, '5'), (2, '3'), (2, '6')]]

*/*Rewrite the following question: What are all the song names by singers who are younger than average?*/*

*/*results:What are all the song names by singers who are younger than average?*/*

*/*Rewrite the following question:Show the citizenship shared by singers with birth year before 1945 and after 1955.*/*

*/*results: Show the citizenship shared by singers with birth year before 1945 and after 1955.*/*

*/*Rewrite the following question:what are the different names of the singers that have sales more than 20000?*/*

*/*results:what are the different names of the singers that have sales more than 20000?*/*

```

/*Given the initial value from Table:*/
['Table TV_Channel: ', [(('700', 'Sky Radio', 'Italy', 'Italian', 'music', '4:3', 'no', 'no', 'Sky Famiglia'),
('701', 'Sky Music', 'Italy', 'Italian', 'music', '4:3', 'no', 'no', 'Sky Famiglia'))]]
['Table TV_series: ', [(1.0, 'A Love of a Lifetime', 'September 24, 2007', '5.8', 9.0, '3.5/9', '9.16', 43.0,
'700'), (2.0, 'Friendly Skies', 'October 1, 2007', '5.3', 9.0, '3.2/8', '8.23', 50.0, '701')]]
['Table Cartoon: ', [(1.0, 'The Rise of the Blue Beetle!', 'Ben Jones', 'Michael Jelenic', 'November14,2008',
101.0, '700'), (2.0, 'Terror on Dinosaur Island!', 'Brandon Vietti', 'Steven Melching', 'November21,2008',
102.0, '701')]]

```

```

/*Rewrite the following question: What are the languages used by the least number of TV
Channels and how many channels use it?*/
/*results:What are the languages used by the least number of TV Channels and how many channels use it?*/

```

/*Here is the input:*/

```

/*Given the initial value from Table:*/
['Table singer: ', [(1, 'Liliane Bettencourt', 1944.0, 30.0, 'France'), (2, 'Christy Walton', 1948.0, 28.8,
'United States'), (3, 'Alice Walton', 1949.0, 26.3, 'United States'), (4, 'Iris Fontbona', 1942.0, 17.4,
'Chile'), (5, 'Jacqueline Mars', 1940.0, 17.8, 'United States')]]

['Table song: ', [(1, 'Do They Know It's Christmas; 1, 1094000.0, 1.0), (2, 'F***k It (I Don't
Want You Back); 1, 552407.0, 1.0), (3, 'Cha Cha Slide', 2, 351421.0, 1.0), (4, 'Call on Me', 4, 335000.0,
1.0), (5, 'Yeah', 2, 300000.0, 1.0)]]

```

Rewrite the following question: What is the count of singers?

results:

A.2 Context-Aware Prompt Generation

We only present examples with DAIL-SQL as the base prompt. For simplicity, only one examples as well as the schema and the content of a single table is listed.

```

/* Some SQL examples are provided based on similar problems: */
/* Answer the following: which us city has the highest population
density */

```

```

SELECT city_name FROM city WHERE population = ( SELECT MAX (
    population ) FROM city );

```

-- ...

```

/* Given the following database schema: */

```

```

CREATE TABLE "countries" (
    "CountryId" INTEGER PRIMARY KEY,
    "CountryName" TEXT,
    "Continent" INTEGER,
    FOREIGN KEY (Continent) REFERENCES continents(ContId)
)

```

-- ...

```

/* With at most 15 example rows in each table: */

```

```

INSERT INTO car_makers (Id, Maker, FullName, Country) VALUES /* 7
    rows omitted */

```

```

(1, 'amc', 'American_Motor_Company', '1'),
(5, 'ford', 'Ford_Motor_Company', '1'),
(4, 'gm', 'General_Motors', '1'),
(15, 'peugeot', 'Peugeot', '3'),
(20, 'triumph', 'Triumph', '7'),
(19, 'toyota', 'Toyota', '4'),
(8, 'nissan', 'Nissan_Motors', '4'),
(12, 'mazda', 'Mazda', '4'),
(16, 'renault', 'Renault', '3'),
(23, 'hyundai', 'Hyundai', '8'),
(18, 'subaru', 'Subaru', '4'),
(21, 'volvo', 'Volvo', '6'),
(6, 'chrysler', 'Chrysler', '1'),
(22, 'kia', 'Kia_Motors', '8'),
(14, 'opel', 'Opel', '2');

-- ...

/* Answer the following: Which model has the highest miles per gallon
   in terms of gasoline consumption? */
SELECT

```

A.3 Execution-Guided Refinement

1. Example prompt for a column-name error

```

/* Received the following exception: no such column: Model
In fact, the following 2 table(s) contain column "model": model_list,
   car_names.
Please analyze the cause. */

```

2. Example prompt for other exceptions.

```

/* Received the following exception: no such function: YEAR
Please analyze the cause. */

```

3. Example prompt for a successful execution

```

/* Here are the execution results (total 5 rows):
('Balmoor', 1)
('Glebe Park', 1)
('Recreation Park', 1)
('Somerset Park', 2)
("Stark's Park", 1)

```

```

Judge whether the results are correct, and you selected the exact
   number of columns asked by the question.
Output Yes if you agree with both, or No otherwise. Follow the answer
   with your analysis of the results, and the meanings of the
   columns you selected and those requested by the question. */

```

4. Prompt asking the LLM to refine the query.

```

/* Revise the SQL query to correct any mistakes: */
SELECT

```

A.4 Additional Experimental Results

We further test our framework with DeepSeek-Coder (Guo et al., 2024) via the OpenAI API⁶. Table 5 and 6 display the results of applying DART-SQL to the base prompt of DAIL-SQL and C3 respectively. Note that while testing on the dev and test splits of Spider with DAIL-SQL, we disable the Question Rewriting and Execution-guided Refinement modules due to the limited rephrasing and refining capabilities of DeepSeek. We keep all modules enabled for other datasets with DAIL-SQL and for all datasets with C3.

Base Prompt: DAIL-SQL	Dev		Test		Realistic		DK	
	EM	EX	EM	EX	EM	EX	EM	EX
baseline	55.4	78.4	52.9	78.8	45.4	68.7	35.0	59.1
DART-SQL	54.9	80.0	53.0	79.2	46.7	69.7	40.0	65.4

Table 5: Results with DeepSeek-Coder using DAIL-SQL as the base prompt. For each metric on a dataset, the highest result is marked bold.

Base Prompt: C3	Dev		Test		Realistic		DK	
	EM	EX	EM	EX	EM	EX	EM	EX
baseline	34.3	73.0	35.3	70.7	36.0	63.6	29.2	57.6
DART-SQL	37.2	74.5	37.0	74.6	40.4	72.6	31.4	67.7

Table 6: Results with DeepSeek-Coder using C3 as the base prompt.

We also execute our framework with GPT-3.5-turbo-16k on the BIRD dataset, which replaces the metric of exact-set-match accuracy (EM) with valid efficiency score (VES). VES measures the execution time of an SQL statement compared to the ground-truth statement while requiring the correctness of the execution results. Results are available in Table 7.

Base Prompt: DAIL-SQL	EX	VES
baseline	43.7	47.0
DART-SQL	45.4	47.6

Table 7: Results with GPT-3.5-turbo-16k and DAIL-SQL on BIRD.

⁶DeepSeek does not provide their own API but instead relies on OpenAI API. When we conducted the experiment, the API service was available at <https://api.deepseek.com/v1> for free. Since the release of DeepSeek-V2, the service requires payment and the url has been changed.